# Learning regression and verification networks for long-term visual tracking

Yunhua Zhang, Dong Wang, Lijun Wang, Jinqing Qi, Huchuan Lu
Dalian University of Technology
Dalian, China

`zhangyunhua@mail.dlut.edu.cn, wdice@dlut.edu.cn, wlj@mail.dlut.edu.cn`
`{jinqin, lhchuan}@dlut.edu.cn`

## Abstract

*Compared with short-term tracking, the long-term tracking task requires determining the tracked object is present or absent, and then estimating the accurate bounding box if present or conducting image-wide re-detection if absent. Until now, few attempts have been done although this task is much closer to designing practical tracking systems. In this work, we propose a novel long-term tracking framework based on deep regression and verification networks. The offline-trained regression model is designed using the object-aware feature fusion and region proposal networks to generate a series of candidates and estimate their similarity scores effectively. The verification network evaluates these candidates and outputs the optimal one as the tracked object with its classification score. It is updated online to adapt to the appearance variations based on newly reliable observations. The similarity and classification scores are combined to obtain the final confidence value, based on which our tracker can determine the absence of the target accurately and conduct image-wide re-detection to capture the target successfully when it reappears. Extensive experiments show that our tracker achieves the best performance on the VOT2018 long-term challenge and state-of-the-art results on the OxUvA long-term dataset.*

## 1. Introduction

Visual tracking is a fundamental problem in computer vision, which has many practical applications including video surveillance, vehicle navigation, to name a few. The tracking algorithms can be roughly divided into two branches: short-term tracking and long-term tracking. For the former one, the tracked object is almost always in the camera's field of view but not necessarily fully visible. The tracking algorithm focuses on estimating the accurate positions and scales of the target in short-term sequences. In recent years, numerous trackers [4, 6, 21, 29, 28, 15] have achieved very promising results in short-term tracking benchmarks (such as OTB [34] and VOT2017 [13]). However, the experi-



Figure 1. Visual results of our tracker in representative frames.

Table 1. Comparisons between recent long-term tracking benchmarks and popular short-term ones.

|  | **LTB35** [19] | **OxUvA** [31] | **OTB100** [34] | **VOT2017** [14] |
|---|---|---|---|---|
| Avg frames | 4196 | 4235 | 590 | 350 |
| Max frames | 29700 | 37440 | 3870 | 1494 |
| Min frames | 1389 | 900 | 72 | 36 |
| absent labels | 12% | 52% | 0% | 0% |
| Avg absent labels | 503 | 2.2 | 0 | 0 |

ments on these benchmarks cannot well evaluate the long-term tracking performance of different trackers, and cannot provide valid references for the realistic tracking systems.

Recently, the importance of long-term tracking has been emphasized and some large-scale datasets have been well constructed to address this issue (such as VOT2018 LTB35 [19] and OxUvA [31]). Besides common challenges in short-term scenarios, in the long-term tracking task, the tracked object requires to be captured in long-term sequences, and also disappears and reappears very frequently. Thus, this task provides more challenges than short-term tracking. Table 1 reports some statistics on the frame length and the number of absent labels in popular and recent short-term (OTB, VOT2017) and long-term (VOT2018 LTB35 and OxUvA) benchmarks. First, the frame length in long-term benchmarks is almost ten times longer than that in short-term benchmarks. In addition, there exist a large amount of *absent* labels in long-term tracking scenarios.

1

Thus, it is vital for long-term trackers to provide a valid confidence score indicating the tracked object is *present* or *absent* and have the capability of image-wide re-detection. Figure 1 provides visual results of our tracker in one representative sequence, which shows that our tracker can effectively identify the tracked object is *present* or *absent* and estimate the accurate bounding box when it is *present*.

Up to now, few works have been done to deal with long-term challenges [18, 20, 7, 10, 22], The LCT [20] and PTAV [7] trackers just track and re-detect the target in a local search region, and cannot re-capture the target successfully after the target moves out of view and reappears again. The CMT [22] and FCLT [18] methods merely exploit a single matching or classification model for the entire tracking process, which makes the tracker easily drift due to online variations. The MUSTer algorithm [10] exploits ensemble models to treat short-term and long-term scenarios separately; however, its performance is not satisfactory mainly due to the adopted hand-crafted features.

To deal with challenges in the long-term scenarios, we attempt to develop a deep-learning-based long-term tracker with an integration of regression and verification networks (Figure 2). The regression network $\mathcal{R}$ learns a generic matching function off-line to robustly handle the common appearance variations of the tracked object during the tracking process. The verification network $\mathcal{V}$ further equips the tracker with a strong discriminative power by online learning. In each frame, $\mathcal{R}$ regresses a series of candidate bounding boxes in a local search region, with their scores measuring the similarities between candidates and the object template. Then, $\mathcal{V}$ learns a classification boundary online to further decide whether the most similar candidate is the true target or a distractor. The final confidence score is outputted by the integration of the scores from both $\mathcal{R}$ and $\mathcal{V}$, and indicates the tracked object is *present* or *absent* in the current frame. This score will be used to invoke the image-wide re-detection scheme if necessary.

To summarize, **the main contributions of this work are presented as follows**: i) A novel long-term tracking framework is developed to combine an offline-learned regression network with an online-updated verification network. The regression model aims to candidate proposal, while the verification one is for target identification. ii) A novel object-specific regression network is proposed to generate a series of candidates being similar with the tracked object, which is offline learned effectively and handles intrinsic appearance variations robustly. iii) A valid confidence score is designed to determine the target is present or not, and to make the tracker dynamically switch between local search and global search. iv) Extensive evaluations on the VOT2018 LTB35 and OxUvA long-term benchmarks demonstrate that our tracker achieves the best performance in comparisons with other competing methods.

## 2. Related Work

**Short-term deep trackers.** Recent deep trackers [21, 6, 4, 29, 28, 15] have achieved promising results in short-term sequences, which are usually categorized into either matching-based [3, 30, 15] or classification-based [32, 33, 21, 6, 4, 29, 28, 5] ones. The former ones attempt to train generalized deep neural networks offline, which find the best candidate being most likely to the object template in each frame. The classification-based trackers learn discriminative correlation filters [6, 4, 29, 28, 5] or CNN-based classifiers online [32, 33, 21] to distinguish the target from the cluttered background. The offline-trained matching models are efficient but not well adapt to online variations. The online-updated classifiers have powerful discriminative abilities but are sensitive to noisy observations. Some works [7, 9] have combined both two networks and achieved high performance in short-term scenarios; however, they cannot work well for long-term tracking [19] due to the limitation of their frameworks and re-detection schemes. In this work, we develop a novel long-term tracking framework to effectively integrate the regression (matching) and verification (classification) networks.

**Long-term tracking.** Until now, few works have been proposed for long-term tracking. The TLD [35] algorithm exploits the 'optical flow'-based tracker for local search and an ensemble of weak classifiers for image-wide re-detection. The MUSTer [10] method utilizes a classifier for short-term localization and key points matching for global search. In addition, the CMT [22] tracker merely conducts key points matching for long-term tracking. The above-mentioned trackers can search the target in the entire image, but their performances are not satisfactory due to the adopted hand-crafted low-level features. The LCT [20] and PTAV [7] trackers are equipped with the re-detection scheme for long-term tracking, but they merely track the targets in a local search region to expect that the lost targets will reappear around the previous location. Thus, these two methods are not able to capture the target any more after it moves out of view. FCLT [18] learns correlation filters online and gradually increases the search range with time, but its performance is still far from state-of-the-art. This work proposes a novel deep-learning-based long-term tracking algorithm, which integrates regression and verification networks effectively. The former one generates a series of candidates based on object-aware feature fusion and region proposal network (RPN) [25][1]. The verification network is an online-learned CNN classier to evaluate the generated candidates. It localizes the tracked object accurately and invokes image-wide re-detection if necessary.

---

[1]RPN is a very popular and effective technique in object detection [25, 17, 24]. SiameseRPN [15] have attempted to exploit the RPN method for short-term tracking and achieved a highly competitive performance. But it needs substantial labeled video data for offline training.
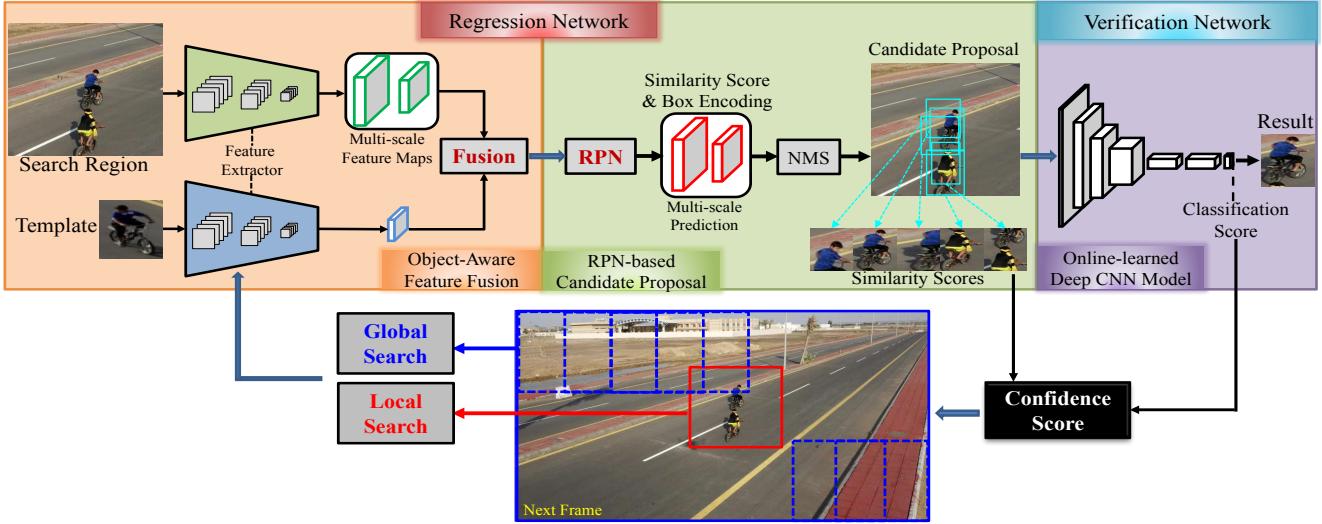
Figure 2. The overall framework of the proposed tracking approach.

# 3. Proposed Tracking Approach

## 3.1. Overview

We propose a long-term tracking framework (Figure 2), with a regression network $\mathcal{R}$ regressing a series of bounding boxes that are similar to the tracked object and a verification network $\mathcal{V}$ verifying the selected candidates.

**Candidate Regression.** The regression network $\mathcal{R}$ is trained off-line to locate a $127 \times 127$ template image within a larger $300 \times 300$ search image. The search region $x$ and the template $z$ are processed by two different transformations $\varphi_1$ and $\varphi_2$, i.e., $\varphi_1(x)$ and $\varphi_2(z)$, respectively. A matching function $f(x, z) = m(\varphi_1(x), \varphi_2(z))$ and a regression function $g(x, z) = r(\varphi_1(x), \varphi_2(z))$ are learned in a single network to densely compare the template image to each candidate region in the search image and regress bounding boxes of those are similar to the target. The function $m$ is a similarity metric and the function $r$ encodes the location information. The most similar ones evaluated by $\mathcal{R}$ are collected to form a candidate pool.

**Verification.** The candidate with the highest similarity score is first cropped out and resized to $107 \times 107$, and verified by the network $\mathcal{V}$. This network learns a classification function online to further filter out distractions appearing during tracking, where $c_i$ indicates the $i$-th candidate. If the most similar one is classified into foreground, the proposed tracker will take it as the tracking result $c^*$ of the current frame. Otherwise, $\mathcal{V}$ selects a foreground candidate from the candidate pool ($[c_1, c_2, ..., c_{N_r}]$, $N_r$ is the number of candidates being considered) with the higher similarity score than other foregrounds as the tracking result.

**Re-detection.** When both $\mathcal{R}$ and $\mathcal{V}$ cannot find any candidate with high similarity and classification scores simulta-

neously, our tracker regards the tracked object being out of view and searches it in the entire image. Unless the tracker has found one patch that is convincing for both $\mathcal{R}$ and $\mathcal{V}$, the tracker regards the target as *absent* in the current frame.

Since $\mathcal{R}$ is fixed (both the parameters and the target template) during the tracking process, it would not accumulate errors and can provide reliable similarity evaluations all the time. $\mathcal{V}$ adapts to variations dynamically, and the inaccurate samples collected online can be regularized by $\mathcal{R}$ to some extent. With the complementation between generalized $\mathcal{R}$ and discriminative $\mathcal{V}$, the proposed tracker is capable of searching the target effectively in long-term sequences. The details are presented in the following sections.

## 3.2. Regression Network

The pipeline of our regression network is shown in Figure 2. It adopts the SSD [17] detection framework and MobileNets [11] as feature extractors. The two streams of the network share the same architecture. Since the sizes of the target are not identical in two inputs (i.e. the search region image patch $x$ and the object template $z$), the two branches use different parameters. The upper branch takes the search region (cropped around the location of the target in the last frame) as input. It outputs two scales of feature maps $\varphi_1(x)$, namely $19 \times 19 \times 512$ and $10 \times 10 \times 512$, to handle drastic scale variations. The lower branch takes the object template (ground truth given in the first frame) as input and outputs a single $1 \times 1 \times 512$ feature vector $\varphi_2(z)$. The feature maps of the object template and the search region are then fused and sent into the region proposal networks (RPNs). The outputs of RPNs are a series of feature maps encoding bounding box information and matching results. Non-maximum-suppression (NMS) is performed afterwards to get the can-

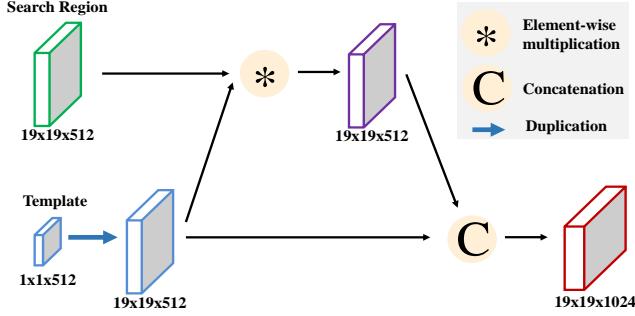didate bounding boxes, with threshold of $IoU$ being 0.6.



Figure 3. Illustration of the fusion procedure.

**Object-Aware Feature Fusion.** The fusion procedure aims to provide the region proposal networks with representative features of the search region and the tracked object conducting similarity measure and bounding box prediction. To merge with the feature maps of the search region, the single $1 \times 1 \times 512$ feature vector of the object template $\varphi_2(z)$ is duplicated to $19 \times 19 \times 512$ and $10 \times 10 \times 512$ feature maps. The obtained feature maps of the object template have the same sizes with these of the search region $\varphi_1(x)$ and are fused with them correspondingly. We take the fusion procedure of the $19 \times 19$ scale for example, as shown in Figure 3. The $19 \times 19 \times 512$ feature maps of the search region are first multiplied with these of the object template, which highlights the locations that are similar to provide similarity information for the latter metric evaluations. The result feature maps are then concatenated with the $19 \times 19 \times 512$ feature maps of the object template to give the latter RPN prior information about the target for accurate bounding box regression. The final $19 \times 19 \times 1024$ feature maps are the inputs of the corresponding RPN. The fusion operation of the $10 \times 10$ scale is the same as that of the $19 \times 19$.

**RPN-based Candidate Proposal.** The proposed tracker has one region proposal network for each scale. Each subnetwork has two branches, one for similarity calculation and the other for proposal regression. The branch for metric learning takes the fused feature maps rather than the original feature maps of the image patch as inputs, which is different from traditional detection frameworks [17, 24]. Each branch consists of three convolutional layers with $3 \times 3$ and $1 \times 1$ kernels. If there are $k$ anchors, for each scale, the network needs to output $2k$ channels for matching and $4k$ channels for regression. When training the network with several anchors and several scales, we employ the loss function that is used in SSD [17]. Our loss function is defined as a sum of the matching loss (cross-entropy loss) and the localization loss (smooth $L_1$ loss with normalized coordinates).

### 3.3. Verification Network

In Section 3.2, the regression network $\mathcal{R}$ generates foreground candidates being similar to the object template. However, the generated candidate pool may contain distractions that lead to drifts. It is inappropriate to directly update $\mathcal{R}$ with respect to a specific object, since the errors are inevitably accumulated as the tracking process proceeds. Thus, a reasonable manner is to keep $\mathcal{R}$ fixed for reliable similarity measure and introduce an additional verification network $\mathcal{V}$ for candidate evaluation. The architecture of $\mathcal{V}$ is similar to that of MDNet [21]. It takes a $107 \times 107$ patch as input and outputs two neurons indicating the probabilities of foreground and background respectively. More details about the network architecture can be found in [21].

Similar with the original MDNet, we update the last three convolutional layers of the network online to train a strong softmax-based classifier which can distinguish the foreground from the background effectively. Through online updating, $\mathcal{V}$ helps the proposed tracker tackle with various cluttered background during tracking. Since the training samples are assessed by both $\mathcal{R}$ and $\mathcal{V}$, $\mathcal{V}$ is not likely to break down due to inappropriate updates.

### 3.4. Tracking Strategy

The proposed tracker first searches the tracked object in the search region, which is four times of the object size. After obtaining the best candidate in each frame, we can consider the tracked object as *present* or *absent* based on its confidence score, and determine how to search the target in the next frame. If the confidence score $S_c$ is below a threshold of $0.3$, the tracker regards the tracked object as *absent* and invokes the global search scheme in the next frame. Otherwise, the tracker treats the tracked object as *present* and continues to adopt the local search in the next frame. As shown in Figure 2, the local and global search schemes are dynamically switched based on the confidence score of the best candidate, which indicates whether the tracker finds a reliable candidate or not.

In this work, the confidence score $S_c$ of the selected candidate in each frame is defined by both the regression score $S_r$ and the verification score $S_v$ as

$$S_c = \begin{cases} 1.0, & S_v > \theta_{v'} \ or \ S_r > \theta_{r'}, S_v > 0 \\ 0, & S_r < \theta_r, S_v < 0 \\ S_r, & otherwise \end{cases}, \quad (1)$$

where $\theta_{v'} = 20.0$, $\theta_{r'} = 0.5$ and $\theta_r = 0.3$. The principles of equation (1) include: (1) When $\mathcal{V}$ is very confident ($S_v > \theta_{v'}$) or both $\mathcal{R}$ and $\mathcal{V}$ are confident enough ($S_r > \theta_{r'}, S_v > 0$), our tracker outputs a confidence score of $1.0$; (2) When both $\mathcal{R}$ and $\mathcal{V}$ give negative feedbacks, the tracker returns a confidence score of $0$; (3) Otherwise, the tracker sets the confidence score from the regression one ($S_c = S_r$).

## 3.5. Implementation Details

**Network Architectures and Pre-trained Parameters.** For $\mathcal{R}$, we use the MobileNet architecture [11] as the feature extractor for both branches. The architecture of $\mathcal{V}$ is VGGM as in [21]. Both VGGM of $\mathcal{V}$ and MobileNets of $\mathcal{R}$ load parameters pretrained on ImageNet classification task [27], while only the regression network $\mathcal{R}$ is further trained offline using external datasets.

**Training Data Preparation.** During the training phase of $\mathcal{R}$, sampled pairs are selected from both ILSVRC [26] image and video object localization datasets with a random interval. For the former dataset, we train $\mathcal{R}$ to make it have the capability to regress any kind of object for a given object template. To be specific, we choose an object of interest from an image randomly and crop the corresponding detection region around the chosen object. For the video object localization dataset, the similarity calculation branch of the regression network $\mathcal{R}$ further learns a generic matching function for tracking to tolerate the common appearance variations. The regression network is trained in an end-to-end manner using the stochastic gradient descent. Because of the need of training the regression branch, some data augmentations are adopted including affine transformation and random erasing [36].

**Hyper Parameters and Training Strategy.** In long-term sequences, since the target moves out of view frequently and its size often changes dramatically when it re-appears, we adopt two scales with different ratios of anchor. The anchor ratios we adopt are $[0.33, 0.5, 1, 2, 3]$. The strategy to pick positive and negative training samples is also important in our proposed framework. The criterion used in object detection task is adopted here that we use $IoU$ together with two thresholds $th_{hi}$ and $th_{lo}$ as the measurement. Positive samples are defined as the anchors which have $IoU > th_{hi}$ with their corresponding ground truth. Negative ones are defined as the anchors which satisfy $IoU < th_{lo}$. We set $th_{lo}$ to 0.5 and $th_{hi}$ to 0.7. Similarly to the training process of SSD [17], instead of using all the negative examples, we sort them using the confidence loss of each default box and pick the top ones so that the ratio between the negatives and positives is at most $3 : 1$. There are totally $500,000$ iterations performed during training phase and the batch size is 32. Each batch consists of 16 pairs from image object localization dataset and 16 pairs from video object localization dataset. We use the $10^{-2}$ and $10^{-3}$ learning rates both for $200,000$ iterations, then continue training for $100,000$ iterations with $10^{-4}$.

**Online Tracking.** During the inference phase, the parameters of the regression network are fixed. The object template for matching is the groundtruth given in the first frame. We only update the verification network in a similar way to that of the original MDNet [21].

## 4. Experiments

Our tracker is implemented using Tensorflow [1] on a PC machine with an Inter i7 CPU (32G RAM) and a NVIDIA TITAN X GPU (12G memory). The average tracking speed is 2.7fps. *We will make our source codes (both training and testing codes) be available to the public.*

We evaluate the proposed method on the VOT-2018 LTB35 dataset [19] and OxUvA Long-term dataset [31]. The detailed comparisons are presented as follows.

### 4.1. State-of-the-art Comparisons on VOT-2018 LTB35 Dataset

The VOT-2018 LTB35 dataset [19] is presented in Visual Object Tracking (VOT) challenge 2018 for evaluating long-term trackers, which includes 35 challenging sequences of various objects (e.g., persons, car, motorcycles, bicycles and animals) with a total frame length of $146847$ frames. Each contains on average 12 long-term target disappearances, and lasts on average 40 frames. Therefore, this dataset can fully evaluates trackers' long-term tracking capabilities, which are essential for tracking in the wild.

Following the evaluation protocol of VOT-2018 [14], the tracking precision (**Pr**), recall (**Re**) and **F-score** metrics are utilized for accuracy evaluation. Based on the precision and recall, the threshold F-measure $F(\tau_\theta)$ is defined as

$$F(\tau_\theta) = 2Pr(\tau_\theta)Re(\tau_\theta)/(Pr(\tau_\theta) + Re(\tau_\theta)), \quad (2)$$

where $\tau_\theta$ is the threshold. Then, the F-score is defined as the highest score on the F-measure plot over all thresholds $\tau_\theta$ (i.e., taken at the tracker-specific optimal threshold). This manner avoids arbitrary manual-set thresholds and then encourages fair evaluation. In VOT-2018 LTB35 [19], the **F-score** is the primary long-term tracking measure and is used for ranking different trackers. In addition, the results of the re-detection experiment are adopted to compare different trackers, including the average number of frames required for re-detection (**Frames**) and the percentage of sequences with successful re-detection (**Success**).

The detailed comparisons are reported in Table 2. From Table 2, we can conclude that our tracker achieves the top-ranked performance in terms of **F-score**, **Pr** and **Re** criteria while maintaining a high re-detection success rate[2]. Especially, the proposed method obtains the best performance among all compared trackers in terms of the **F-score** measure, which is the most important metric on the VOT-2018 LTB35 dataset [19]. The DaSiam_LT method achieves a slightly worse **F-score** than our tracker (0.607 for DaSiam_LT and 0.610 for ours), however, the success rate of DaSiam_LT is almost zero, which means the DaSiam_LT tracker fails in the re-detection experiment. In contrast, the

---

[2]The **Re** values of DaSiam_LT and Ours are 0.588216 and 0.587921, respectively. Thus, we mark Ours by red and DaSiam_LT by blue.

5

Table 2. Performance evaluation for 15 state-of-the-art algorithms on the VOT-2018 LTB35 dataset [19]. The best three results are marked in **red**, **blue** and **green** bold fonts respectively. The trackers are ranked from top to bottom using the **F-score** measure.

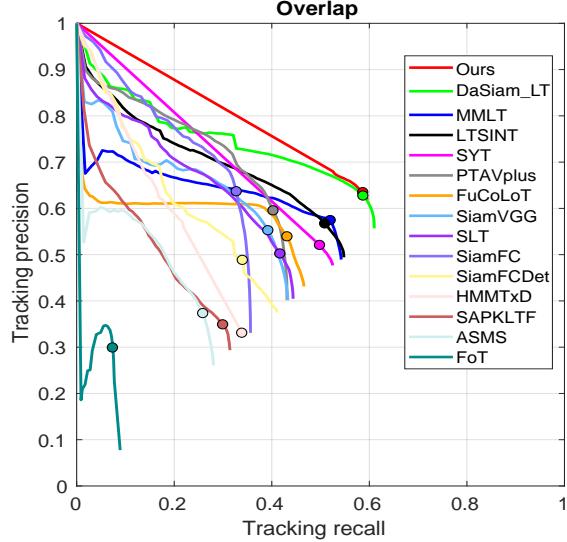| Tracker | F-score | Pr | Re | Frames (Success) |
|---|---|---|---|---|
| **Ours** | **0.610** | **0.634** | **0.588** | 1 (100%) |
| DaSiam_LT | **0.607** | **0.627** | **0.588** | - (0%) |
| MMLT | **0.546** | 0.574 | **0.521** | 0 (100%) |
| LTSINT | 0.536 | 0.566 | 0.510 | 2 (100%) |
| SYT | 0.509 | 0.520 | 0.499 | 0 (43%) |
| PTAVplus | 0.481 | 0.595 | 0.404 | 0 (11%) |
| FuCoLoT | 0.480 | 0.539 | 0.432 | 78 (97%) |
| SiamVGG | 0.459 | 0.552 | 0.393 | - (0%) |
| SLT | 0.456 | 0.502 | 0.417 | 0 (100%) |
| SiamFC | 0.433 | **0.636** | 0.328 | - (0%) |
| SiamFCDet | 0.401 | 0.488 | 0.341 | 0 (83%) |
| HMMTxD | 0.335 | 0.330 | 0.339 | 3 (91%) |
| SAPKLTF | 0.323 | 0.348 | 0.300 | - (0%) |
| ASMS | 0.306 | 0.373 | 0.259 | - (0%) |
| FoT | 0.119 | 0.298 | 0.074 | - (6%) |



Figure 4. Average precision-recall curves of our trackers and other state-of-the-art methods on VOT-2018 LTB35 [19]. Different trackers are ranked based on the maximum of the F-Score.

proposed method achieves a success rate of 100%. We also note that our tracker merely uses the ImageNet Detection [26] and ILSVRC [26] video datasets for training, while the DaSiam_LT method adopts much more data for training including ImageNet Detection [26], ILSVRC [26], Youtube-BB [23] and COCO Detection [16].

To visualize the superiority of our tracker, we provide representative results of our tracker and other two top-ranked methods reported in [14], along with their confidence scores on challenging image sequences. As shown in Figure 5, the targets disappear frequently in these videos. Our tracker outputs a very low confidence score when the target is absent, while other trackers still give relatively higher confidence scores and drift to distractions (such as *ballet*, *skiing*, *following*, *freestyle* and *cat1*). In sequence *group1*, our method captures the tracked object successfully but the other two trackers drift to other people easily. Thus, the confidence scores outputted by our tracker is able to indicate the presence of the tracked object accurately.

## 4.2. Ablation Study

In this subsection, we conduct ablation analysis to evaluate different components of our tracker. With different experiment settings, we design four variants of our tracker, which are respectively named as "Ours w/o Verification", "Ours w/o Concatenation", "Ours w/o Multiplication" and "Ours (Siamese)". The meanings of these notions are explained as follows. (1) "Ours w/o Verification" denotes the tracker that only uses the regression network for tracking, which means that the tracked object is localized as the bounding box with the highest confidence score pre-

dicted by the regression network in each frame. (2) "Ours w/o Concatenation" and "Ours w/o Multiplication" stand for different operations used in the fusion procedure of the regression network. The former one represents the tracker that only uses the result of multiplication between the feature of the search region and that of the template without concatenating the feature of the template. "Ours w/o Multiplication" denotes the tracker that only concatenates the feature of the search region with that of the template without multiplication. (3) "Ours (Siamese)" indicates the tracker in which the feature extractors for the object template and search region share the same parameters.

Table 3. Ablation analysis of the proposed tracker on the VOT-2018 LTB35 dataset. The best results are marked in the **red** font.

| Tracker | F-score | Pr | Re |
|---|---|---|---|
| Ours | **0.610** | **0.634** | **0.588** |
| Ours w/o Verification | 0.525 | 0.563 | 0.493 |
| Ours w/o Concatenation | 0.582 | 0.630 | 0.540 |
| Ours w/o Multiplication | 0.442 | 0.504 | 0.394 |
| Ours (Siamese) | 0.497 | 0.533 | 0.486 |

The results of different variants on VOT-2018 LTB35 [19] are presented in Table 3, from which we can obtain the following conclusions. (1) "Ours w/o Verification" achieves an **F-score** of 0.525, which is also competitive compared with other state-of-the-art methods (shown in Table 2). It means that our regression network can provide not bad metric evaluation. In addition, comparing "Ours" with "Ours w/o Verification", we can conclude that the proposed verification model greatly improves the performance of our
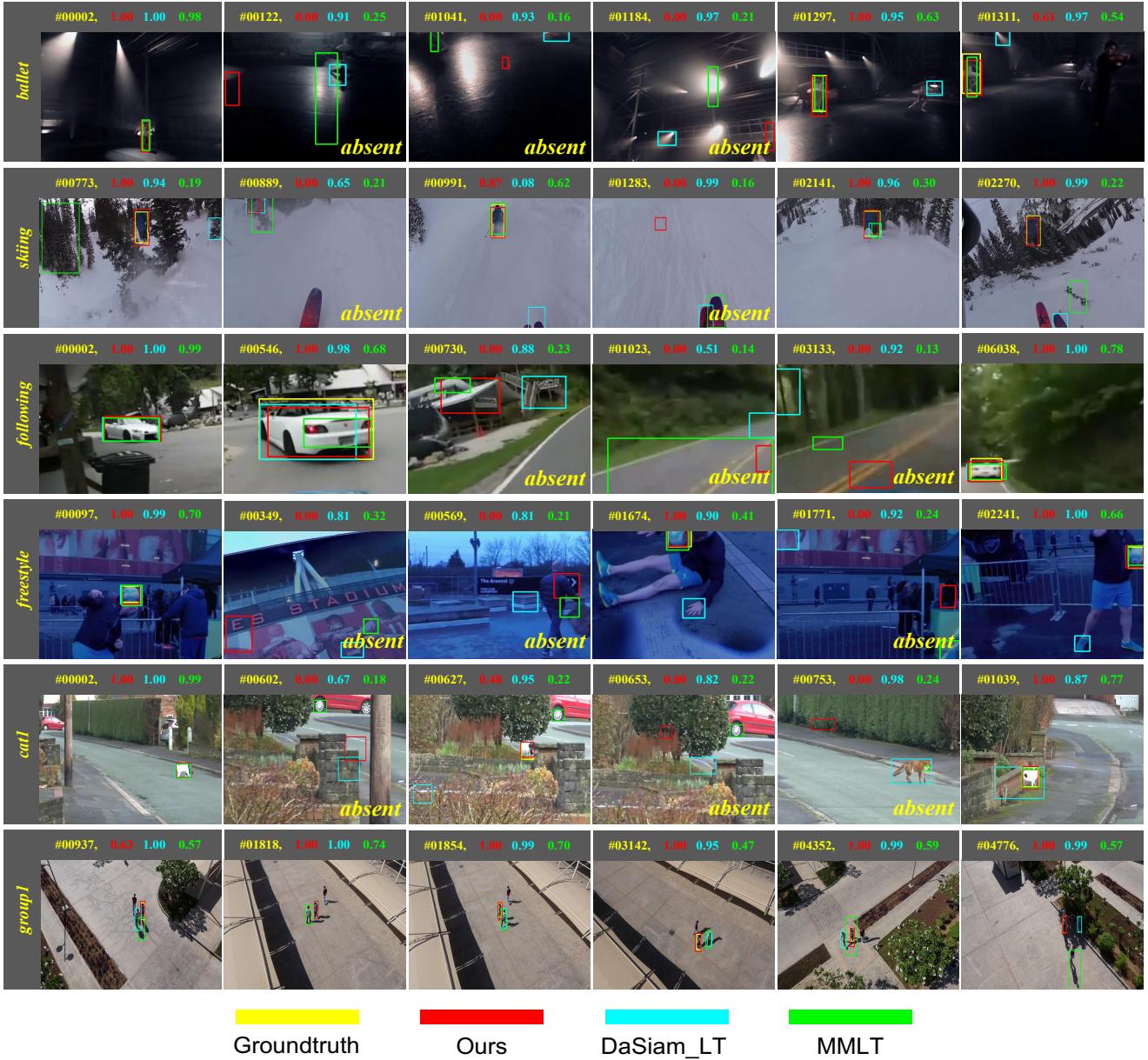
Figure 5. Qualitative results of our tracker, along with DaSiam_LT and MMLT (the top-ranked trackers reported in [14]) on six challenging sequences. From top to bottom: *ballet*, *skiing*, *following*, *freestyle*, *cat1*, *group1*. The confidence score of each tracker is shown in the top of each image. We use "absent" to denote frames in which the target is absent. Best viewed in color with zoom-in for more details.

tracker for long-term scenario. (2) Comparing "Ours" with "Ours w/o Concatenation" and "Ours w/o Multiplication", we can conclude that both concatenation and multiplication operations are essential in the feature fusion module. (3) The comparison between "Ours" and "Ours (Siamese)" shows that the Siamese architecture (i.e., the feature extractors of the template and the search region share the same parameters) leads to inferior performance, mainly caused by different input sizes of these two branches.

## 4.3. State-of-the-art Comparisons on OxUvA Long-term Dataset

The OxUvA long-term dataset [31] comprises 366 object tracks in 337 videos selected from YoutubeBB [23] dataset. The tracklets are labelled at a frequency of 1Hz which is the same as YoutubeBB [23]. In this dataset, each video lasts for average 2.4 minutes, seven times longer than OTB-100 [34] (the most popular short-term tracking dataset). Besides the challenging factors in short-term tracking, severe

Table 4. Detailed comparisons of different trackers on the Ox-UvA [31] dataset. The best three results are marked in **<span style="color:red">red</span>**, **<span style="color:blue">blue</span>** and **<span style="color:green">green</span>** bold fonts respectively. The trackers are ranked from top to bottom using the **MaxGM** measure.

| Tracker | MaxGM | TPR | TNR |
|---------|-------|-----|-----|
| **Ours** | **0.544** | **0.609** | **0.485** |
| SiamFC+R | **0.454** | **0.427** | 0.481 |
| TLD | **0.431** | 0.208 | **0.895** |
| LCT | 0.396 | 0.292 | **0.537** |
| MDNet | 0.343 | **0.472** | 0 |
| SINT | 0.326 | 0.426 | 0 |
| ECO-HC | 0.314 | 0.395 | 0 |
| SiamFC | 0.313 | 0.391 | 0 |
| EBT | 0.283 | 0.321 | 0 |
| BACF | 0.281 | 0.316 | 0 |
| Staple | 0.261 | 0.273 | 0 |



Figure 6. Accuracy plots of different trackers on OxUvA [31].

out-of-view and full occlusion introduce extra challenges which are close to practitioners' needs. The OxUvA dataset is split into *dev* and *test* sets with 200 and 166 tracks respectively. Using these subsets, two challenges have been defined: constrained and open. For constrained challenge, trackers can be developed using only data from OxUvA *dev* set (long-term videos). For the open challenge, trackers can use any public dataset for training except for the YoutubeBB *validation* set, from which OxUvA is constructed. Since all the trackers reported in OxUvA [31] are tested for the open challenge, we also conduct experiment on the open challenge for fair comparison. Following the evaluation benchmark of OxUvA [31], we adopt the true positive rate (**TPR**), true negative rate (**TNR**), and maximum geometric mean (**MaxGM**) to compare different trackers. **TPR** measures the fraction of *present* objects that are reported as *present*, and **TNR** gives the fraction of *absent* objects that are reported as *absent*. The **MaxGM** is defined as,

$$\mathbf{MaxGM} = \max_{0 \leq p \leq 1} \sqrt{((1-p) \cdot \mathbf{TPR})((1-p) \cdot \mathbf{TNR} + p)} , \tag{3}$$

which provides a more informative comparison of different trackers and is used to rank them. A larger **MaxGM** value means a better performance.

We compare our tracker with ten competing algorithms reported in [31], including LCT [20], EBT [37], TLD [35], ECO-HC [4], BACF [8], Staple [2], MDNet [21], SINT [30], SiamFC [3] and SiamFC+R [31]. Among these trackers, LCT [20], EBT [37] and TLD [35] have different re-detection schemes for long-term tracking. ECO-HC [4], BACF [8] and Staple [2] are three recent short-term correlation filter trackers with high accuracy and efficiency. MDNet [21], SINT [30] and SiamFC [3] are three
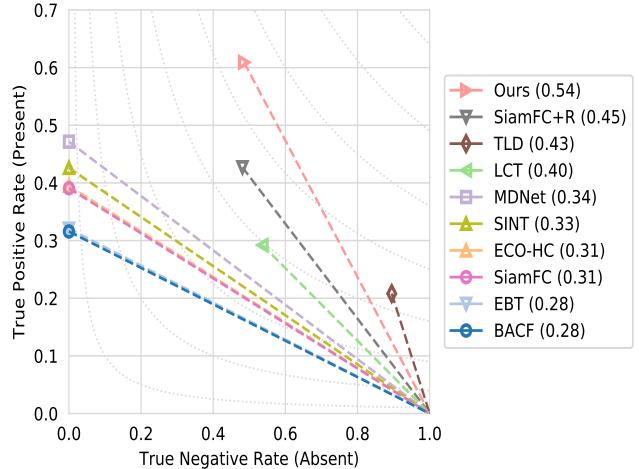
popular algorithms based on deep convolutional networks. SiamFC+R [31] is implemented by equipping SiamFC with a sample re-detection scheme similar to [12]. If the maximum score of the SiamFC's response is smaller than a threshold, the tracker enters the object absent mode. In this situation, SiamFC+R attempts to find the tracked object within a random search region in each frame until the maximum score again surpasses the threshold, at which point the tracker returns to the object present mode.

The detailed comparisons are reported in Figure 6 and Table 4. We can see that our tracker achieves the top-ranked performance in terms of **MaxGM** and **TPR** while maintaining a competitive **TNR**. Especially, our tracker has the best performance among all compared trackers in terms of **MaxGM**, which is the most important metric on the OxUvA dataset. Compared with the second best method (SiamFC+R), our tracker achieves a substantial improvement, with a relative gain of $19.82\%$ in **MaxGM**.

## 5. Conclusion

In this paper, we propose a novel long-term tracking framework based on regression and verification networks. The regression network is trained offline and fixed online for effective candidate proposal and reliable similarity estimation without accumulating errors in long sequences. In addition, the verification network is exploited to evaluate the generated candidates and fine-tuned online to capture the appearance variations. A dynamic switch scheme between local search and image-wide re-detection is designed based on a confidence score that is determined by the outputs from both regression and verification networks. The experimental results on two recent long-term benchmarks demonstrate that our tracker achieves significantly better performance than other state-of-the-art algorithms, which can be acted as a new baseline for further studies.

# References

[1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, and M. Devin. Tensorflow: Large scale machine learning on heterogeneous distributed systems. In *arXiv preprint arXiv:1603.04467*, 2016.

[2] L. Bertinetto, J. Valmadre, S. Golodetz, O. Miksik, and P. H. S. Torr. Staple: Complementary learners for real-time tracking. In *CVPR*, 2016.

[3] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. S. Torr. Fully-convolutional siamese networks for object tracking. In *ECCV Workshop*, 2016.

[4] M. Danelljan, G. Bhat, F. S. Khan, and M. Felsberg. ECO: Efficient convolution operators for tracking. In *CVPR*, 2017.

[5] M. Danelljan, G. Hger, F. S. Khan, and M. Felsberg. Convolutional features for correlation filter based visual tracking. In *ICCV Workshop*, 2015.

[6] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *ECCV*, 2016.

[7] H. Fan and H. Ling. Parallel tracking and verifying: A framework for real-time and high accuracy visual tracking. In *ICCV*, 2017.

[8] H. K. Galoogahi, A. Fagg, and S. Lucey. Learning background-aware correlation filters for visual tracking. In *ICCV*, 2017.

[9] A. He, C. Luo, X. Tian, and W. Zeng. A twofold siamese network for real-time object tracking. In *CVPR*, 2018.

[10] Z. Hong, Z. Chen, C. Wang, X. Mei, D. Prokhorov, and D. Tao. Multi-store tracker (muster): A cognitive psychology inspired approach to object tracking. In *CVPR*, 2015.

[11] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. In *arXiv preprint arXiv:1704.04861*, 2017.

[12] J. S. S. III and D. Ramanan. Tracking as online decision-making: Learning a policy from streaming videos with reinforcement learning. In *ICCV*, 2017.

[13] M. Kristan, A. Leonardis, J. Matas, M. Felsberg, R. Pflugfelder, L. Zajc, T. Vojir, and G. Hager. The visual object tracking vot2017 challenge results. In *ICCV Workshop*, 2017.

[14] M. Kristan, A. Leonardis, J. Matas, M. Felsberg, R. Pflugfelder, L. Zajc, T. Vojir, and G. Hager. The sixth visual object tracking vot2018 challenge results. In *ECCV Workshop*, 2018.

[15] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu. High performance visual tracking with siamese region proposal network. In *CVPR*, 2018.

[16] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.

[17] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *ECCV*, 2016.

[18] A. Lukei, L. . Zajc, T. Voj, J. Matas, and M. Kristan. FCLT - a fully-correlational long-term tracker. In *arXiv preprint arXiv:1711.09594*, 2017.

[19] A. Lukei, L. . Zajc, T. Voj, J. Matas, and M. Kristan. Now you see me: evaluating performance in long-term visual tracking. In *arXiv preprint arXiv:1804.07056*, 2018.

[20] C. Ma, X. Yang, C. Zhang, and M. H. Yang. Long-term correlation tracking. In *CVPR*, 2015.

[21] H. Nam and B. Han. Learning multi–domain convolutional neural networks for visual tracking. In *CVPR*, 2016.

[22] G. Nebehay and R. Pflugfelder. Clustering of static-adaptive correspondences for deformable object tracking. In *CVPR*, 2015.

[23] E. Real, J. Shlens, S. Mazzocchi, X. Pan, and V. Vanhoucke. Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. In *CVPR*, 2017.

[24] J. Redmon and A. Farhadi. YOLO9000: Better, faster, stronger. In *CVPR*, 2017.

[25] S. Ren, K. He, R. Girshick, and J. Sun. Faster R–CNN: towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[26] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, and M. Bernstein. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2014.

[27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, and M. Bernstein. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.

[28] C. Sun, D. Wang, H. Lu, and M. H. Yang. Correlation tracking via joint discrimination and reliability learning. In *CVPR*, 2018.

[29] C. Sun, D. Wang, H. Lu, and M. H. Yang. Learning spatial-aware regressions for visual tracking. In *CVPR*, 2018.

[30] R. Tao, E. Gavves, and A. W. M. Smeulders. Siamese instance search for tracking. In *CVPR*, 2016.

[31] J. Valmadre, L. Bertinetto, J. F. Henriques, R. Tao, A. Vedaldi, A. W. Smeulders, P. H. Torr, and E. Gavves. Long-term tracking in the wild: a benchmark. In *ECCV*, 2018.

[32] L. Wang, W. Ouyang, X. Wang, and H. Lu. Visual tracking with fully convolutional networks. In *ICCV*, 2015.

[33] L. Wang, W. Ouyang, X. Wang, and H. Lu. Stct: Sequentially training convolutional networks for visual tracking. In *CVPR*, 2016.

[34] Y. Wu, J. Lim, and M. H. Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1834–1848, 2015.

[35] K. M. Zdenek Kalal and J. Matas. Tracking-learning-detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1409–1422, 2012.

[36] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang. Random erasing data augmentation. In *arXiv preprint arXiv:1708.04896*, 2017.

[37] G. Zhu, F. Porikli, and H. Li. Beyond local search: Tracking objects everywhere with instance-specific proposals. In *CVPR*, 2016.